

# QUALITY PROGRAMMABLE VECTOR PROCESSORS FOR APPROXIMATE COMPUTING

**Swagath Vekataramani**<sup>1</sup>, Vinay Chippa<sup>1</sup>, Srimat  
Chakradhar<sup>2</sup>, Kaushik Roy<sup>1</sup>, Anand Raghunathan<sup>1</sup>

<sup>1</sup>**Integrated Systems Laboratory  
School of ECE, Purdue University**

<sup>2</sup>**NEC Laboratories America**

**International Symposium on Microarchitecture 2013**



# COMPUTERS == PRECISE CALCULATORS

TASK

$$\frac{433}{21} > 20.4$$



```
float x = 433/21  
float y = 20.4  
(x > y) ? YES :NO
```



YES

$$\frac{433}{21} > 1$$



```
float x = 433/21  
float y = 1  
(x > y) ? YES :NO
```



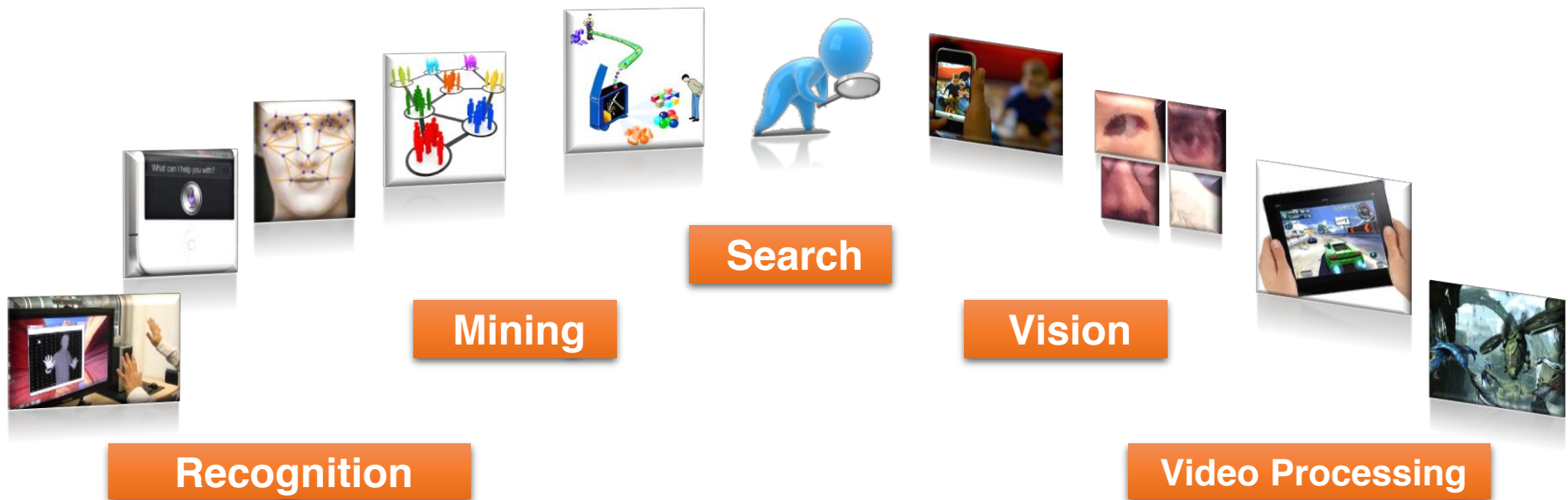
YES

But, I worked  
**harder** than  
needed



- ▶ Leads to inefficiency
- ▶ And, an **overkill** (for many applications)

# EVOLVING APPLICATION LANDSCAPE



- ▶ Relaxed notion of correctness
  - Results cannot be arbitrary either

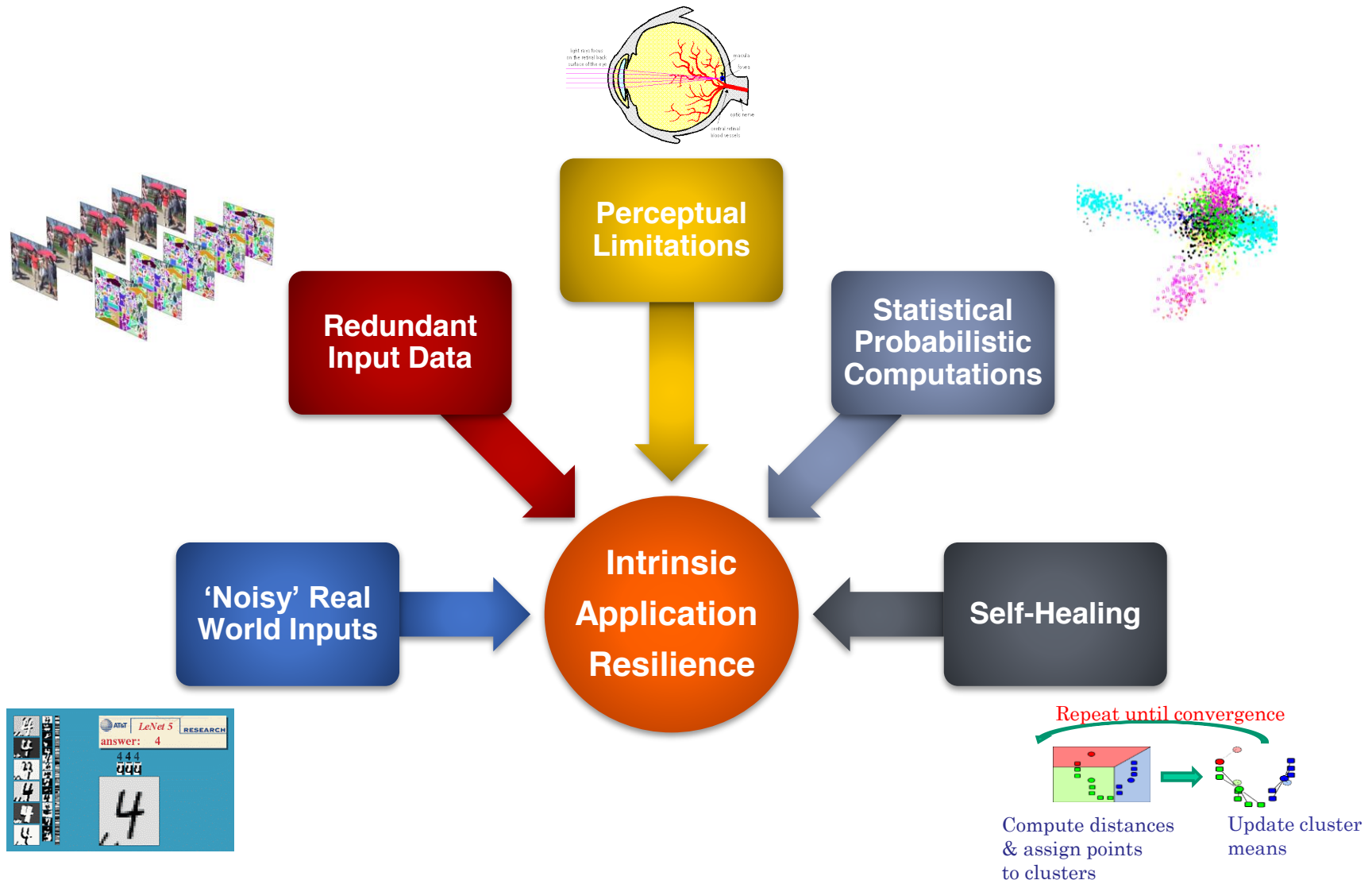
*Good enough* answers !!!

# INTRINSIC APPLICATION RESILIENCE: A NEW DIMENSION TO OPTIMIZE HW & SW



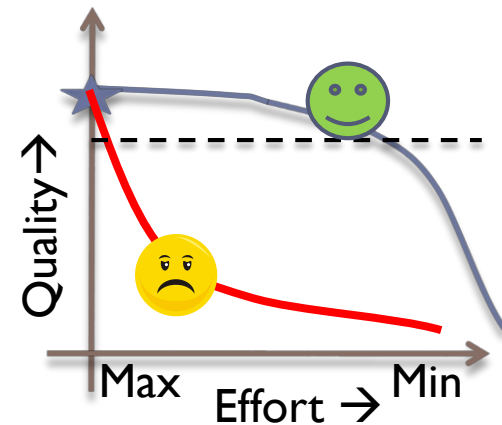
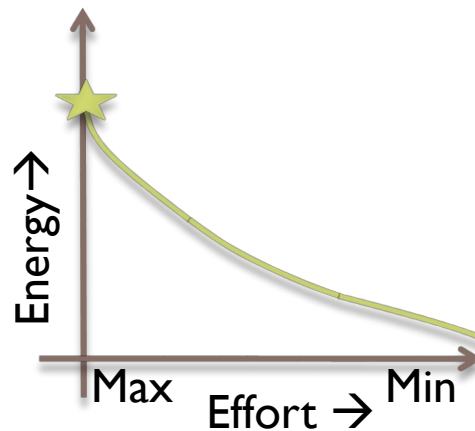
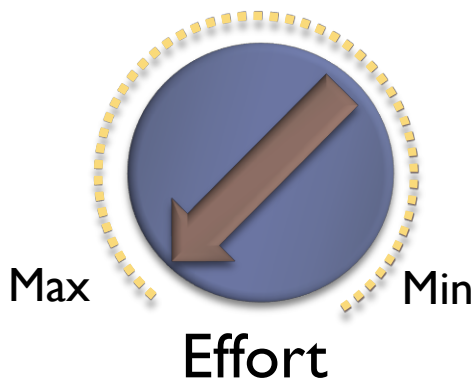
- ▶ The ability to produce outputs of acceptable quality despite many of their computations executed *imprecisely*

# INTRINSIC APPLICATION RESILIENCE: SOURCES

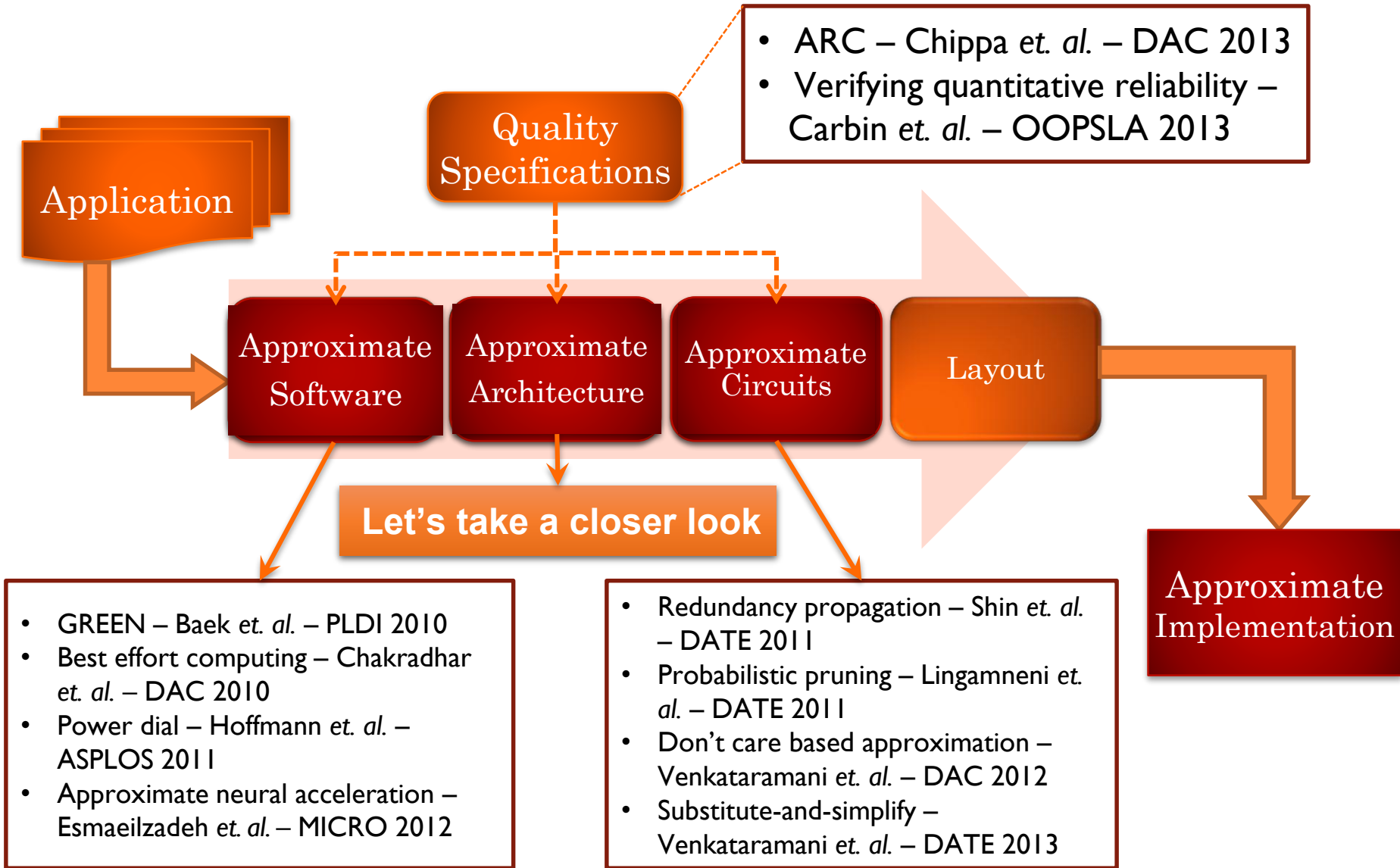


# APPROXIMATE COMPUTING: DESIGN PHILOSOPHY

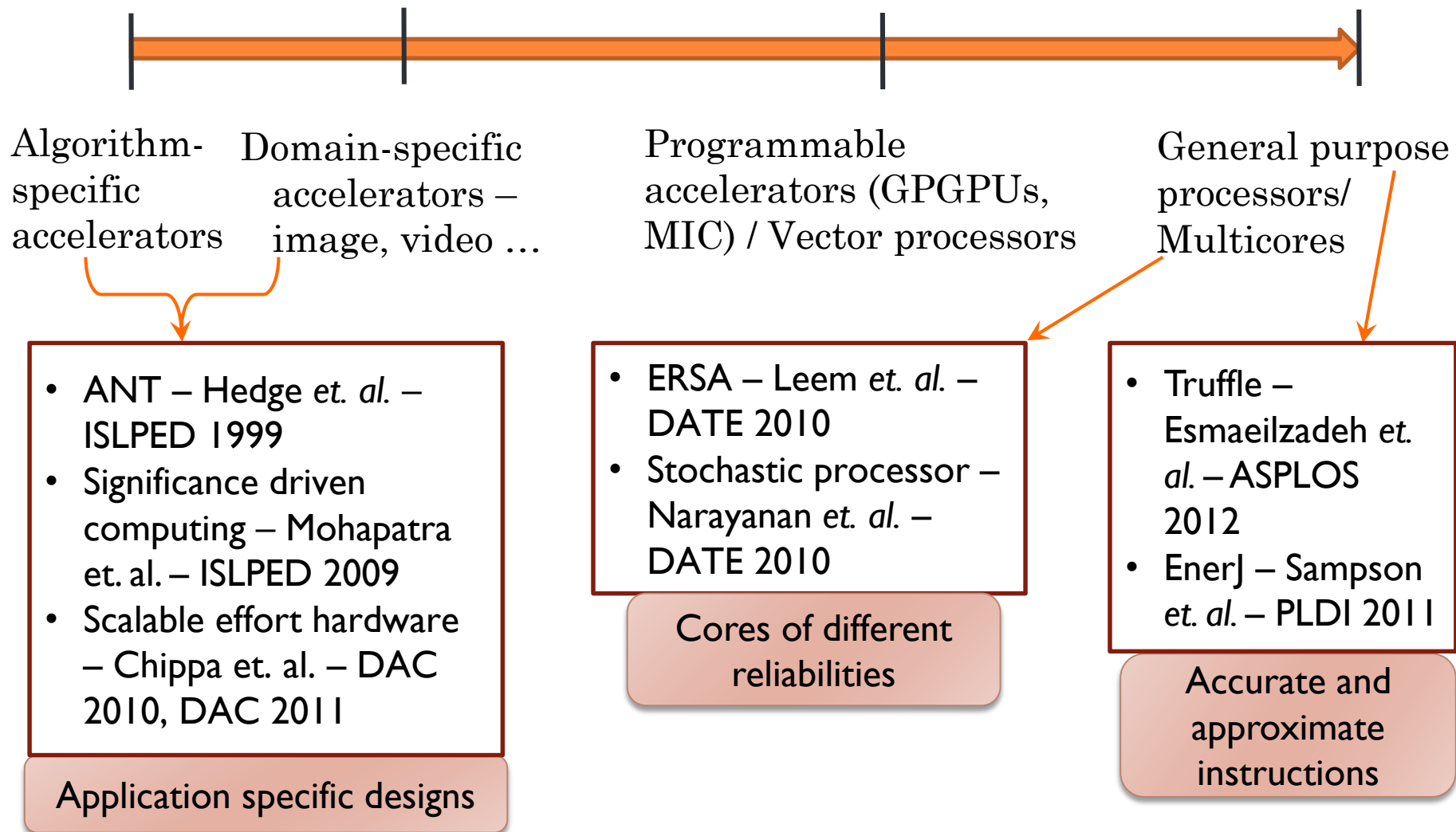
- ▶ Systems that can modulate the effort expended towards quality of results
  - Higher effort  $\rightarrow$  Higher quality but higher energy
- ▶ How do we get the best Q vs. E tradeoff?
  - Disproportionate benefit



# APPROXIMATE COMPUTING DESIGN TECHNIQUES: OVERVIEW

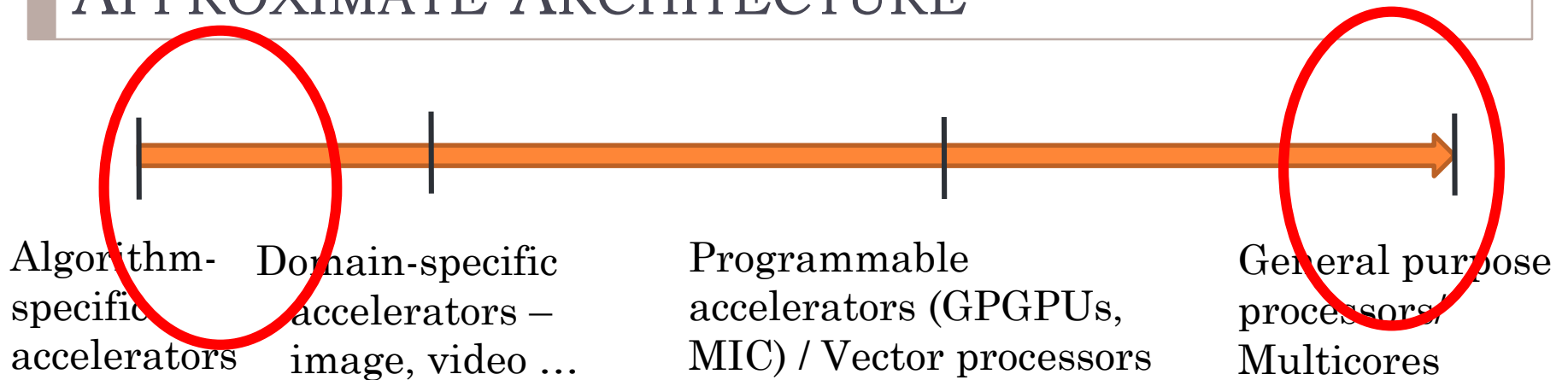


# APPROXIMATE ARCHITECTURE





# APPROXIMATE ARCHITECTURE



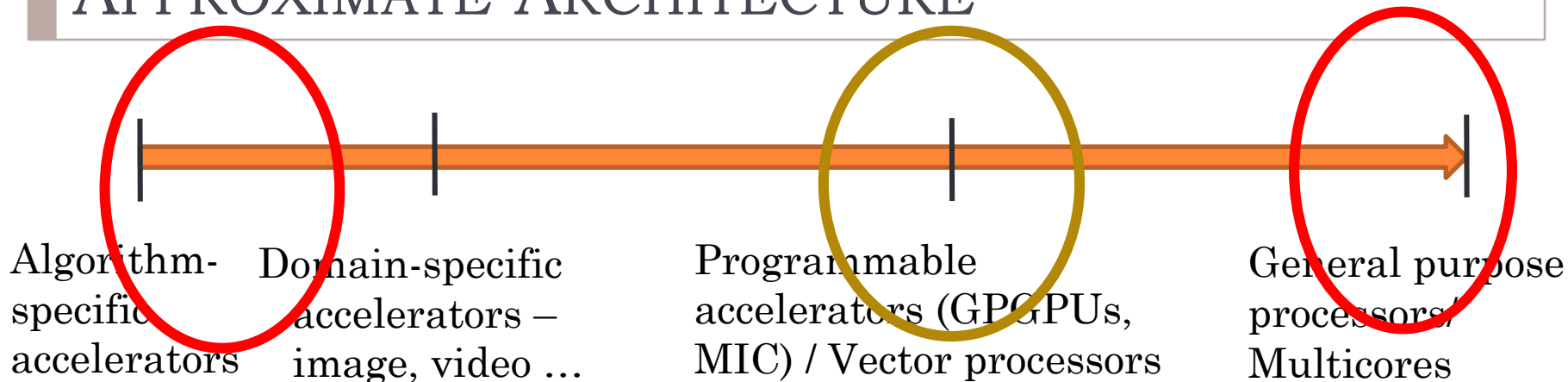
## Pros:

- ☺ Large energy benefits
- ☺ Broader applicability

## Challenges:

- ☹ Limited applicability
- ☹ Inherently limited energy benefit – Dominated by **control front-ends** that cannot be approximated
- ☹ Allow arbitrary errors in hardware – **limits** the fraction of computations that can be approximated

# APPROXIMATE ARCHITECTURE



## Opportunity:

☺ Wide range of applications – fine grained parallelism

☺ **SIMD**: Control overheads amortized over many execution units

☺ Need **quality guarantees** from HW  
– We will address that!

# CONTRIBUTIONS

- ▶ **Quality programmable processors**

An abstract model for programmable approximate processors

- ▶ **QUORA**

A quality programmable 1D/2D vector processor

# Quality programmable processors

## Requirements:

- HW/SW interface for applications to expose resilience
- Micro-architecture that can translate resilience to efficiency



# QUALITY PROGRAMMABLE PROCESSORS

- ▶ **Quality Programmability:** Ability to specify beyond what can be accurate & approximate to HW
- ▶ Notion of *quality* explicitly built into the instruction set

HW/SW INTERFACE

## Quality Programmable ISA

- Quality fields in instructions

qpADD dest, op1, op2, **MAG, 1%**

Quality-programmable  
add

Error magnitude < 1%  
of maximum  
numeric value of output

- Purely based on instruction semantics

# NEED FOR QUALITY PROGRAMMABILITY

- ▶ Errors injected randomly in x86 instructions – arbitrary vs. bounded

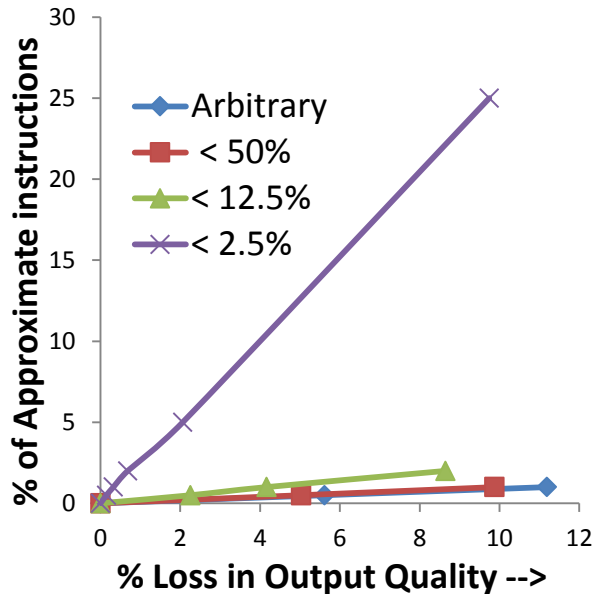
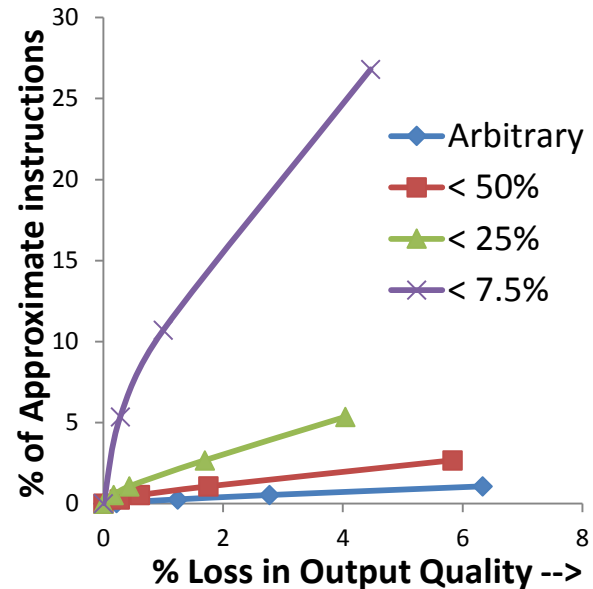


Image Segmentation  
(K-means)



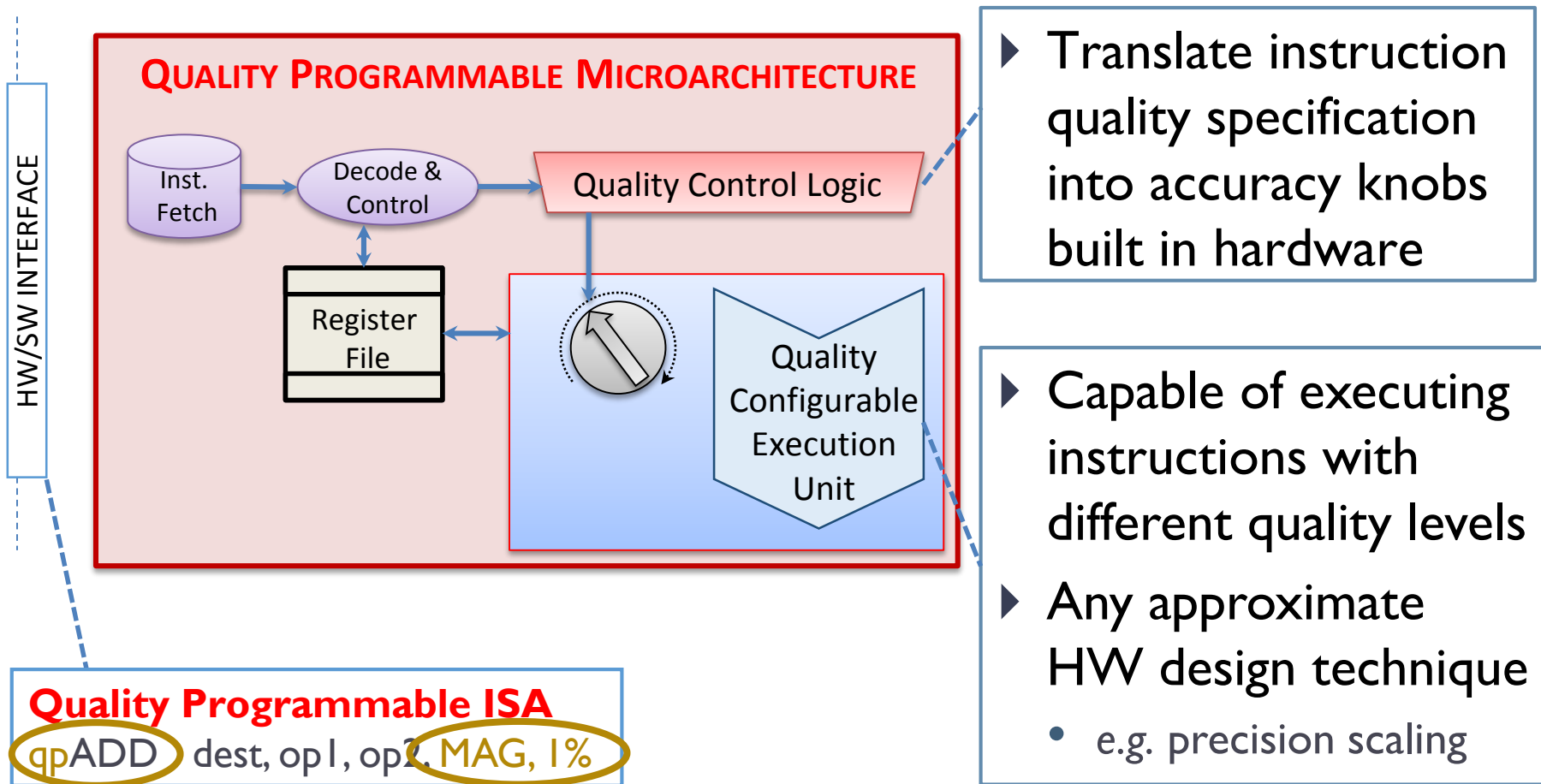
Handwritten Digit Recognition  
(SVM)

25-100X improvement in number of approximate instructions

Constraining errors → much greater opportunity to approximate!!!

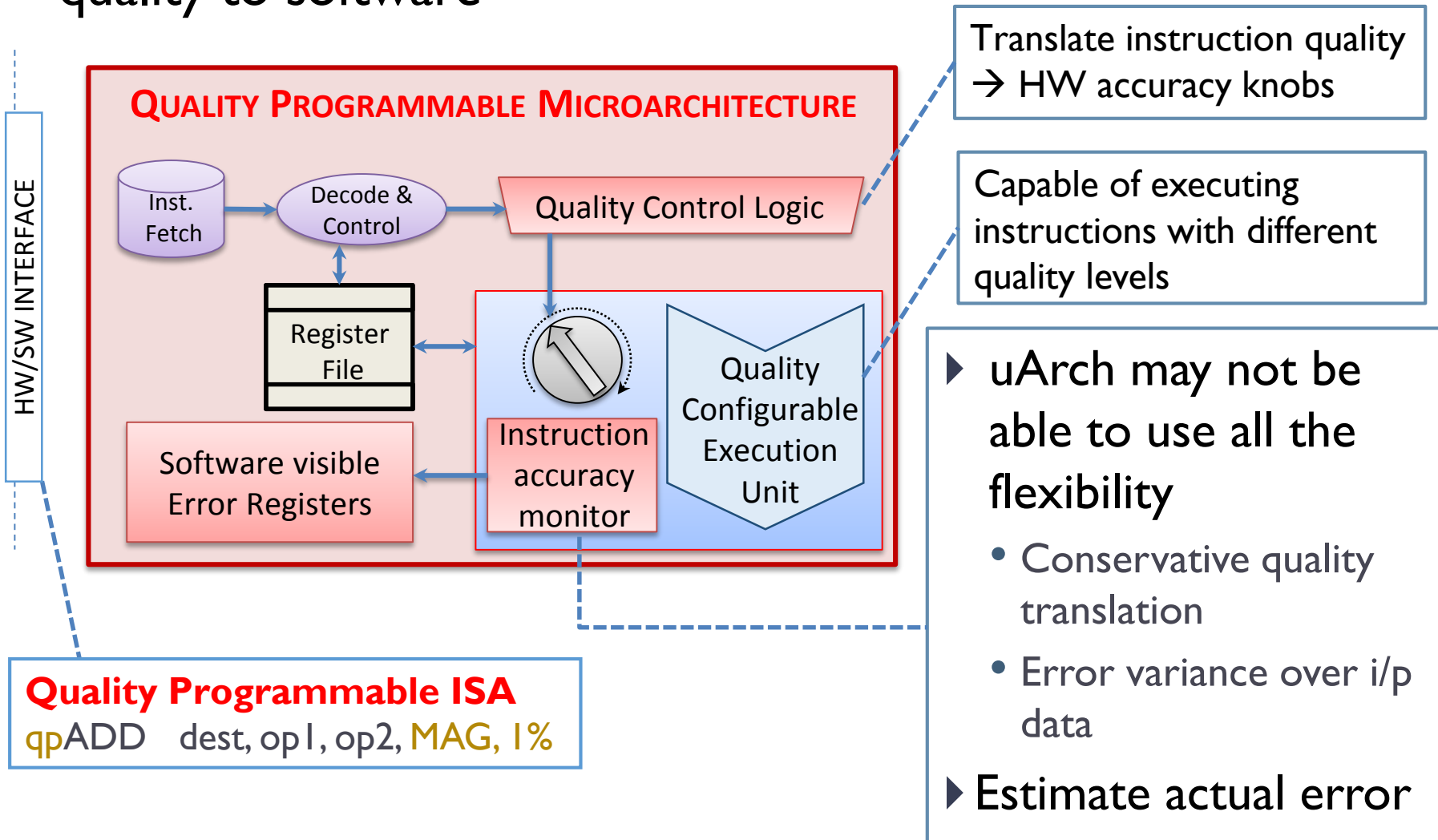
# QUALITY PROGRAMMABLE MICRO-ARCHITECTURE

- ▶ Micro-architecture *guarantees* instruction-level quality



# QUALITY MONITORS AND ERROR FEEDBACK

- ▶ Micro-architecture provides *feedback* on instruction-level quality to software







# QUORA

Quality programmable 1D/2D vector processor

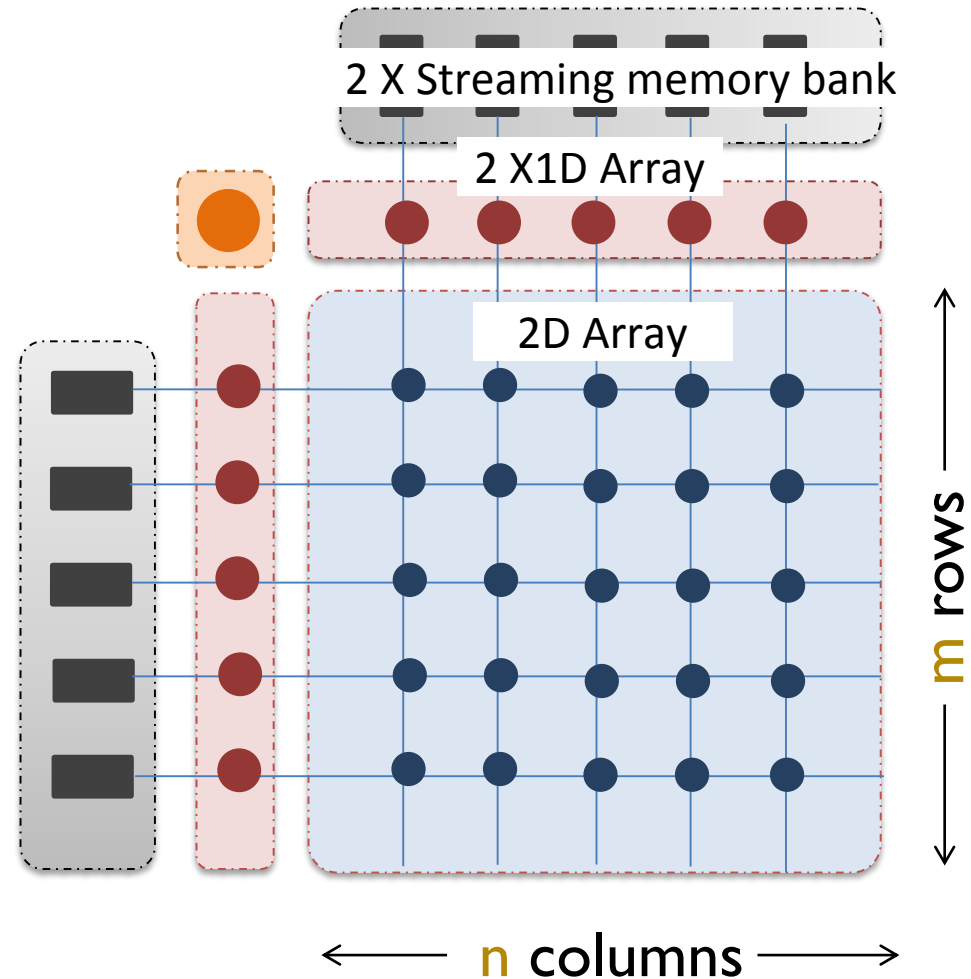


# QUORA: OVERVIEW

## ▶ 3-tier processing element hierarchy

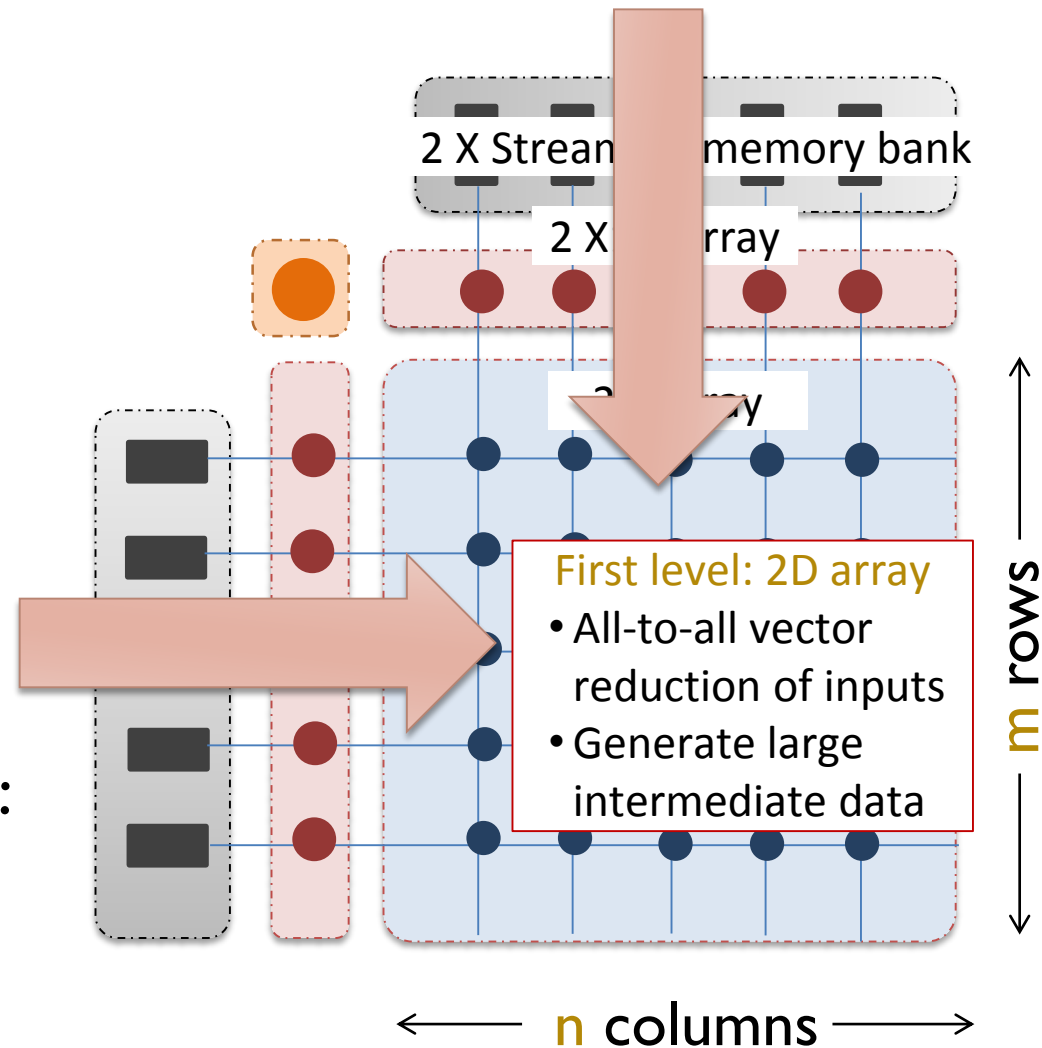
- 2D array PEs
- 2 sets of 1D array PEs
- One scalar PE

## ▶ 2 streaming memory banks along the array borders



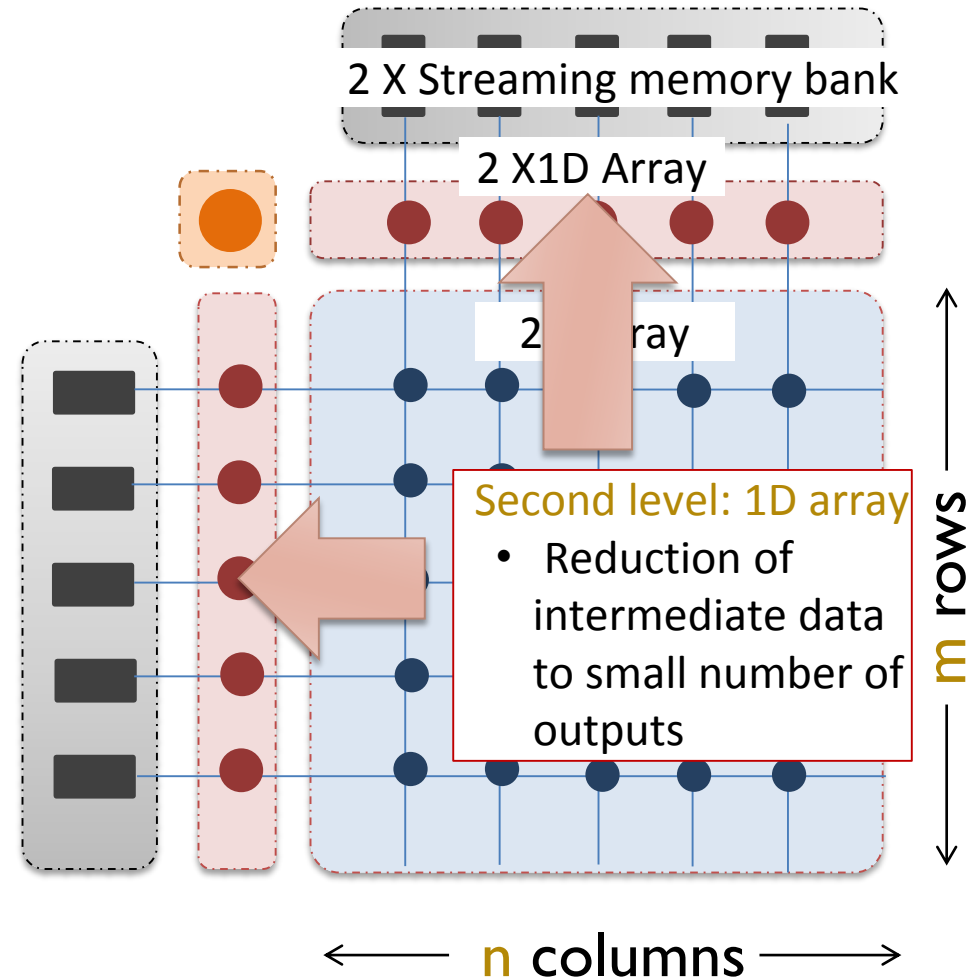
# QUORA: OVERVIEW

- ▶ 3-tier processing element hierarchy
  - 2D array PEs
  - 2 sets of 1D array PEs
  - One scalar PE
- ▶ 2 streaming memory banks along the array borders
- ▶ **Application characteristic:** 2 levels of reduction operations

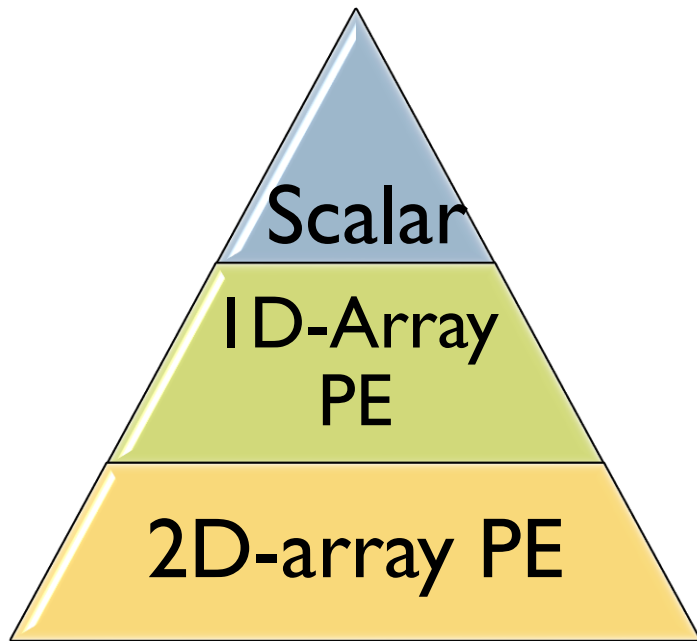


# QUORA: OVERVIEW

- ▶ 3-tier processing element hierarchy
  - 2D array PEs
  - 2 sets of 1D array PEs
  - One scalar PE
- ▶ 2 streaming memory banks along the array borders
- ▶ **Application characteristic:** 2 levels of reduction operations



# PROCESSING ELEMENT HIERARCHY



Functionality / Size



Similar to scalar  
uProcessor

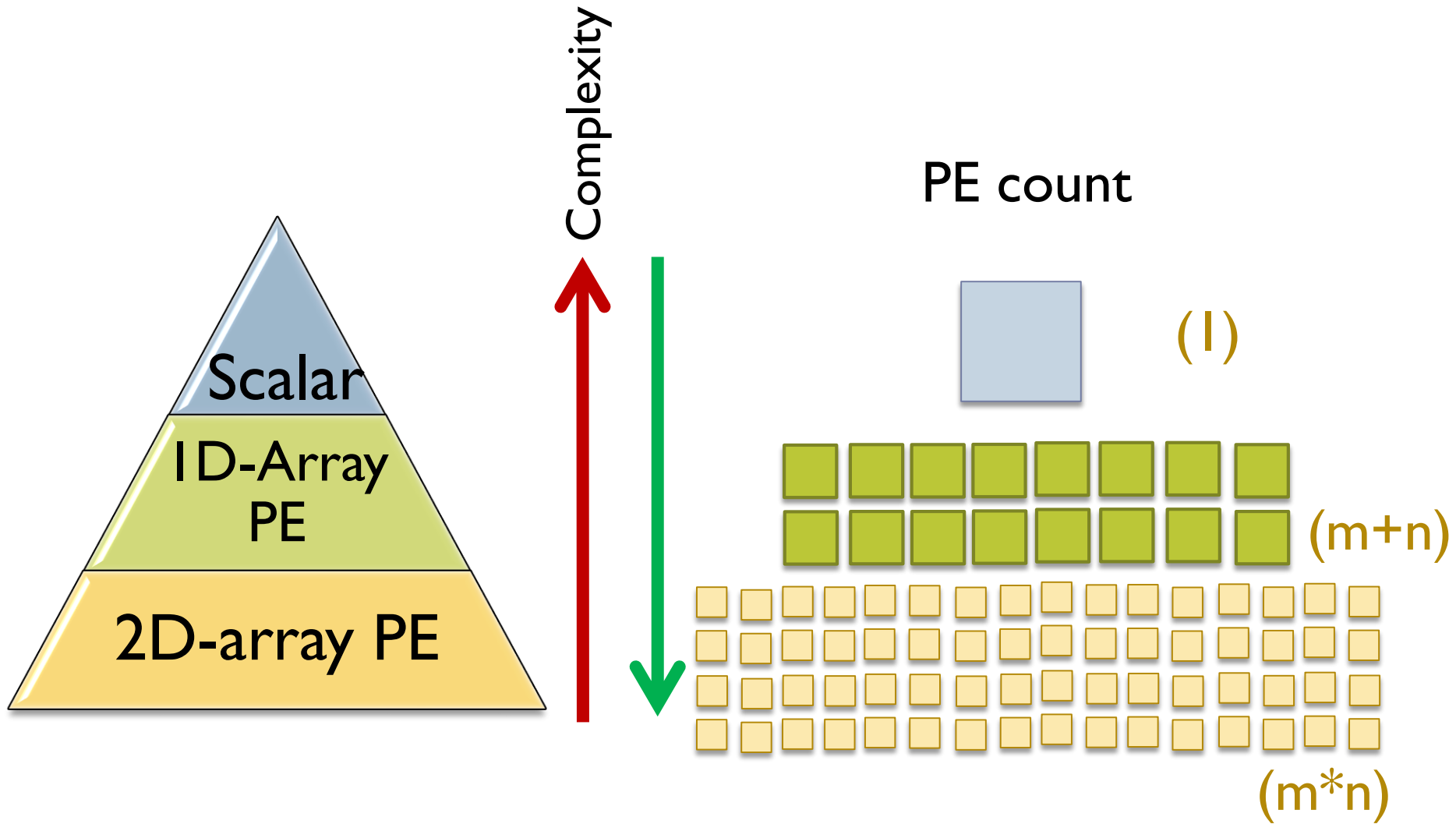


Small register  
file, complex  
execution units

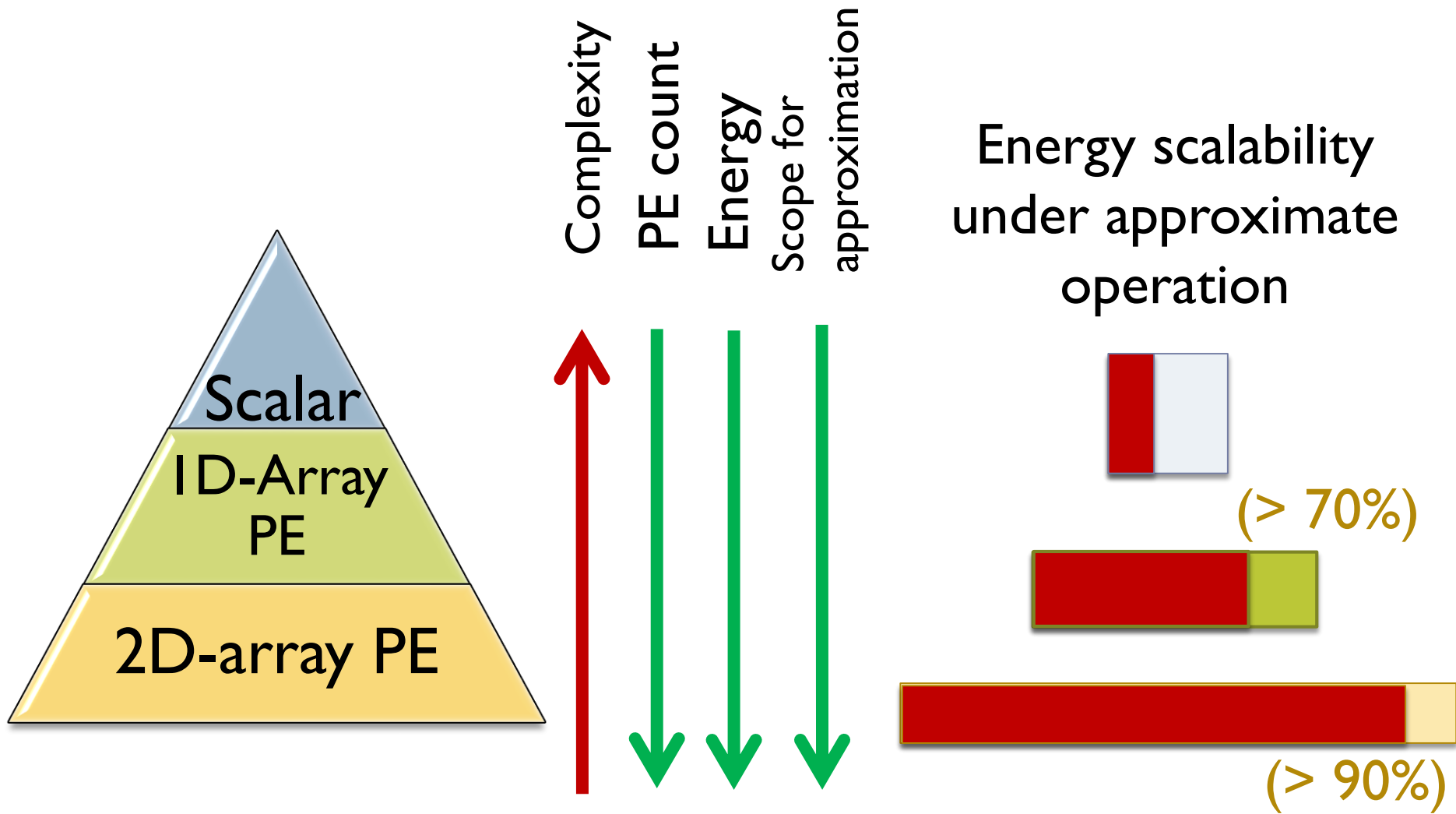


Simple  
accumulator  
based data path

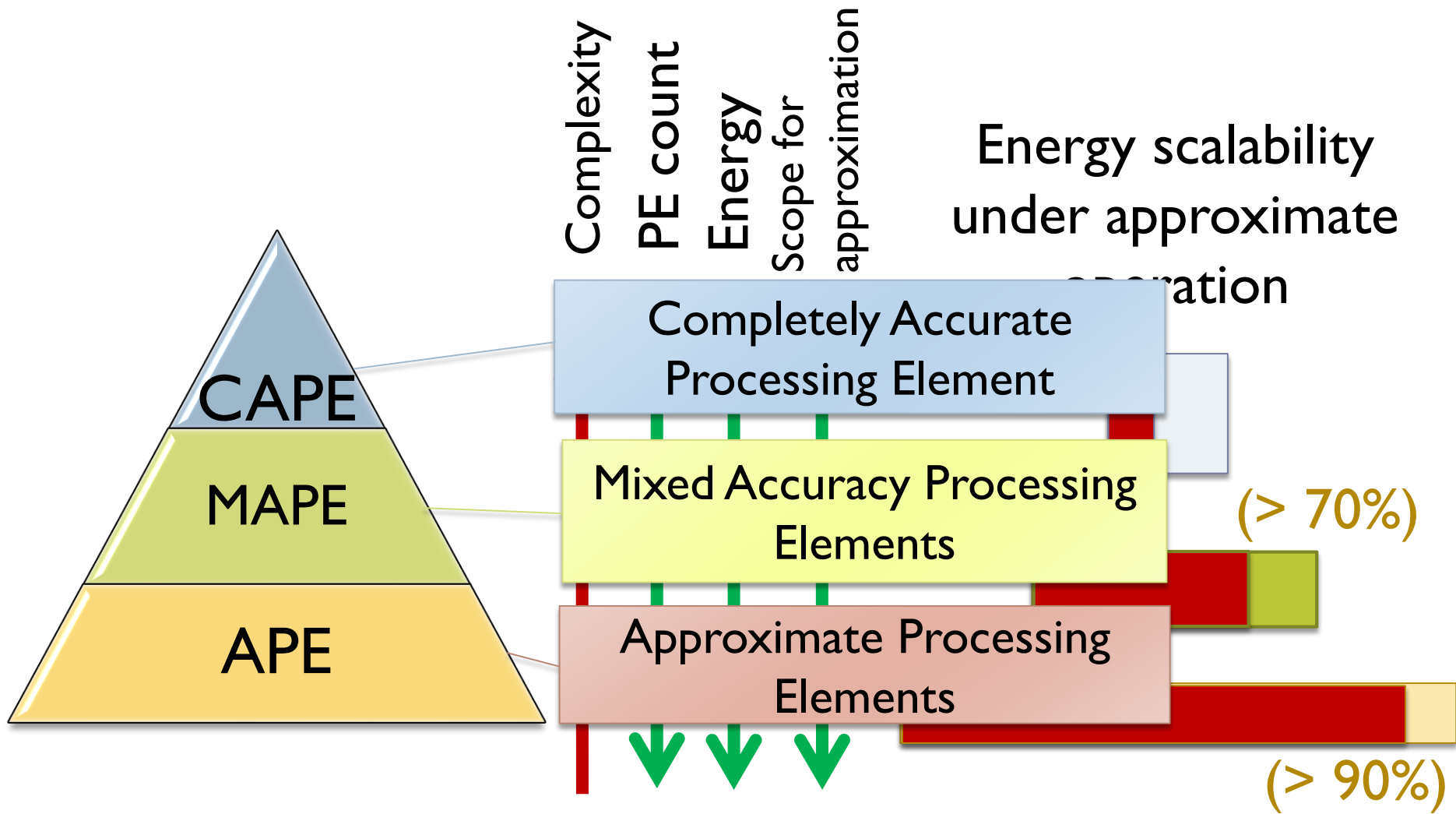
# PROCESSING ELEMENT HIERARCHY



# PROCESSING ELEMENT HIERARCHY



# PROCESSING ELEMENT HIERARCHY



3-tiered PE hierarchy enables larger energy benefits from approximate computing (while matching application characteristics)



# QUORA: INSTRUCTION SET ARCHITECTURE

► **47** Instructions – **9** APE, **22** MAPE, **13** CAPE, **3** SM

Inst. Type	Instruction	Inst. Type	Instruction
Scalar Instructions	LDRI Rd, value	ID Array Reduction Instructions	<b>qpACC</b> <r/c>, R_row_enb, R_col_enb, <b>R_q_type, R_q_amt</b>
	ADDR Rd, Rs1, Rs2		<b>qpMIN</b> <r/c>, R_row_enb, R_col_enb, <b>R_q_type, R_q_amt</b>
	BEZ Rs, Rel. address		
	HALT		
Streaming Memory instructions	LDSM R_length, stride, burst, R_st_add	ID Array Streaming Instructions	SEQ R_length, SReg, R_row_enb, R_col_enb
2D Array Instructions	<b>qpMAC</b> R_length, R_row_enb, R_col_enb, <b>R_q_type, R_q_amt</b>	ID Array Self-Operand Instructions	MVASR <r/c>, R_<r/c>_enb, SReg
	<b>qpMOD2</b> R_length, R_row_enb, R_col_enb, <b>R_q_type, R_q_amt</b>		<b>qpADDX</b> <r/c>, R_<r/c>_enb, Sreg, <b>R_q_type, R_q_amt</b>
	STR <r/c>, R_stride, R_burst, R_st_add, R_row_enb, R_col_enb		<b>qpMUL</b> <r/c>, R_<r/c>_enb, Sreg, <b>R_q_type, R_q_amt</b>
			STMCG <r/c>, R_<r/c>_enb, SReg

# QUORA – QUALITY PROGRAMMABLE INSTRUCTIONS

- ▶ APE and MAPE instructions extended with 2 additional quality fields

e.g. **qpMAC** R\_length, R\_row\_enb, R\_col\_enb, **R\_q\_type, R\_q\_amt**

- Type of error – 3 quality metrics

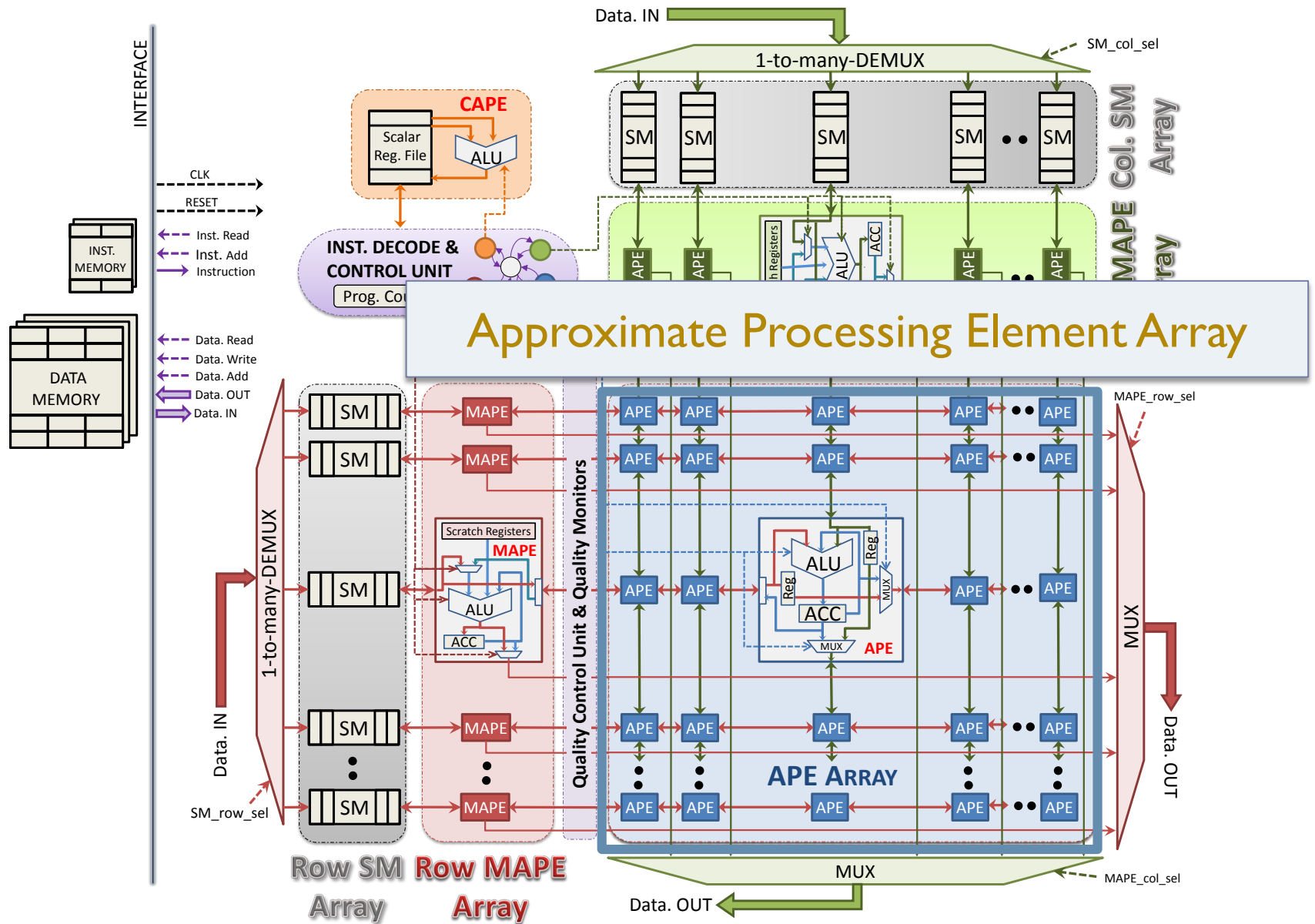
$$MaxErr = |O_{orig} - O_{approx}|$$

$$AveErr = \frac{\sum_{\forall inp} |O_{orig} - O_{approx}|}{Total\ no.\ of\ inputs}$$

$$ErrProb = \frac{No.\ of\ inputs\ for\ which\ O_{orig} \neq O_{approx}}{Total\ no.\ of\ inputs}$$

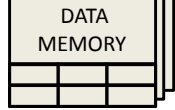
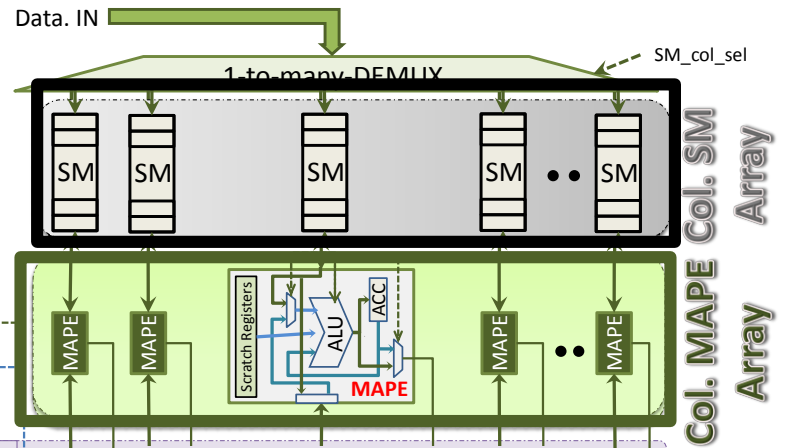
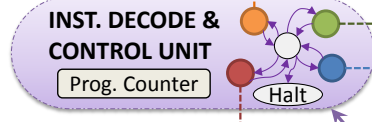
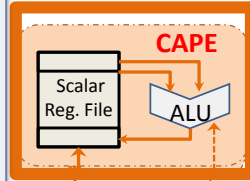
- Amount of error

# QUORA: QP 1D/2D VECTOR PROCESSOR

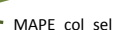
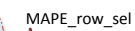
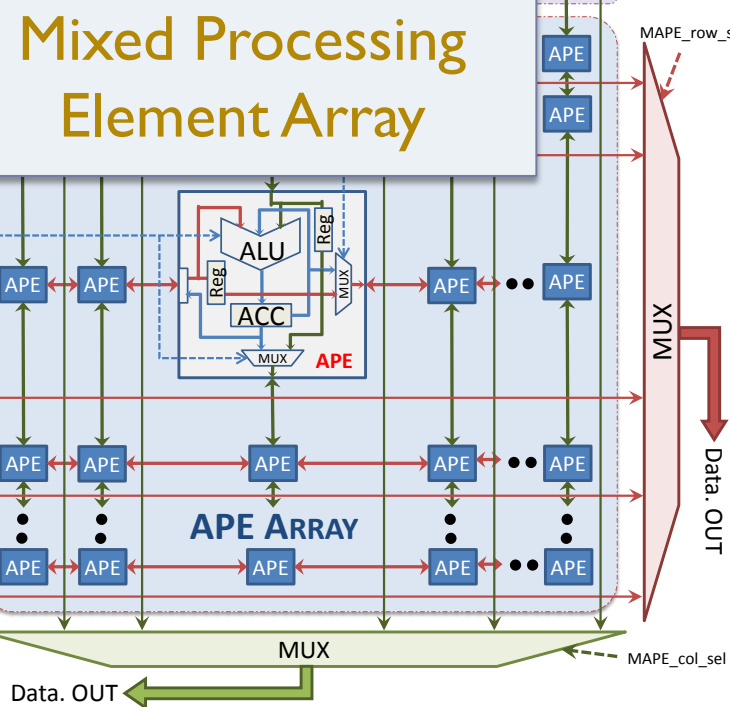
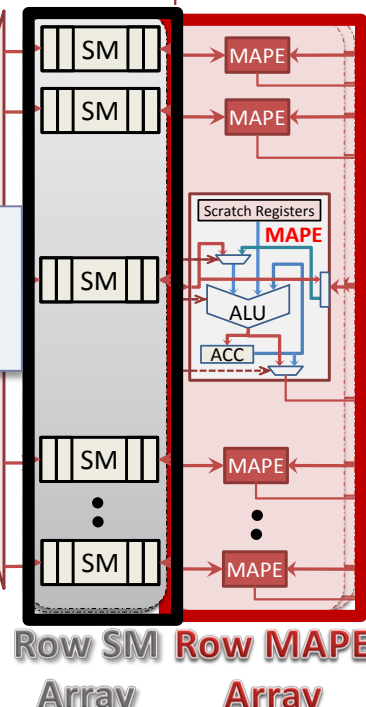
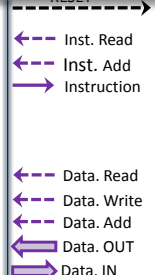


# QUORA: QP 1D/2D VECTOR PROCESSOR

Completely Accurate Processing Element



Streaming Memory Banks



Quality Control & Quality Monitor

Col. SM Array

Col. MAPE Array

Row SM Array

Row MAPE Array

Mixed Processing Element Array

APE ARRAY

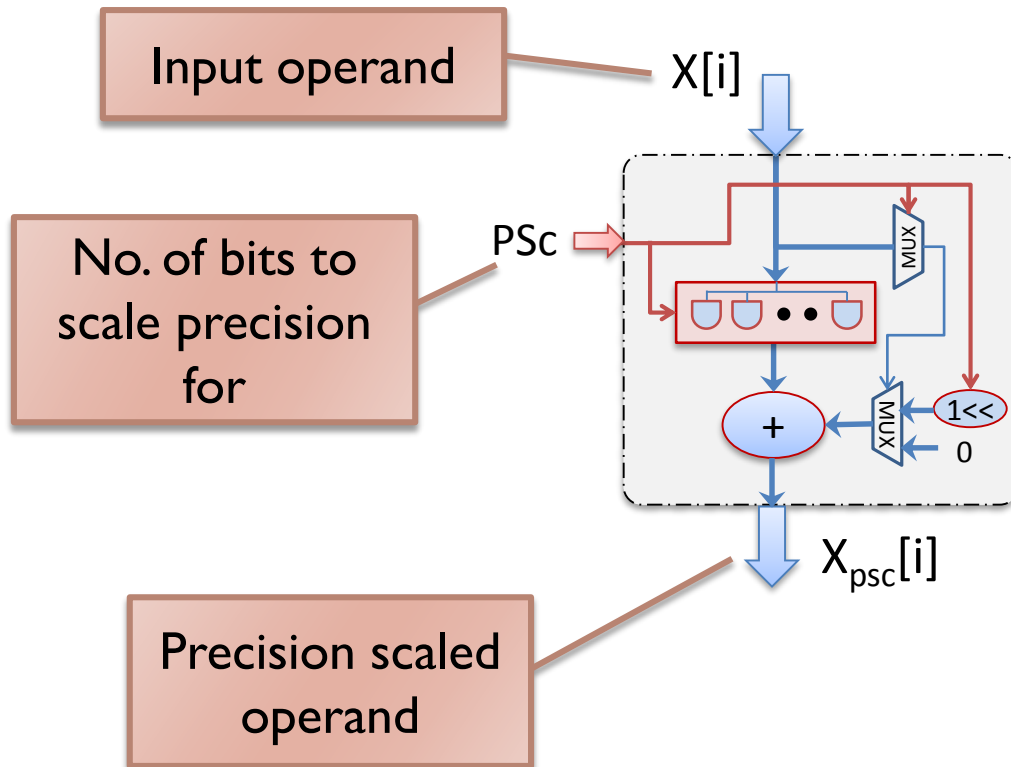
MUX

MAPE\_col\_sel



# QUALITY CONFIGURABLE EXECUTION: PRECISION SCALING

- ▶ Scale the precision of input operands to APEs/MAPEs
  - 4 different flavors
- ▶ Up/down operand round-off



Up/Down Precision  
Scaling

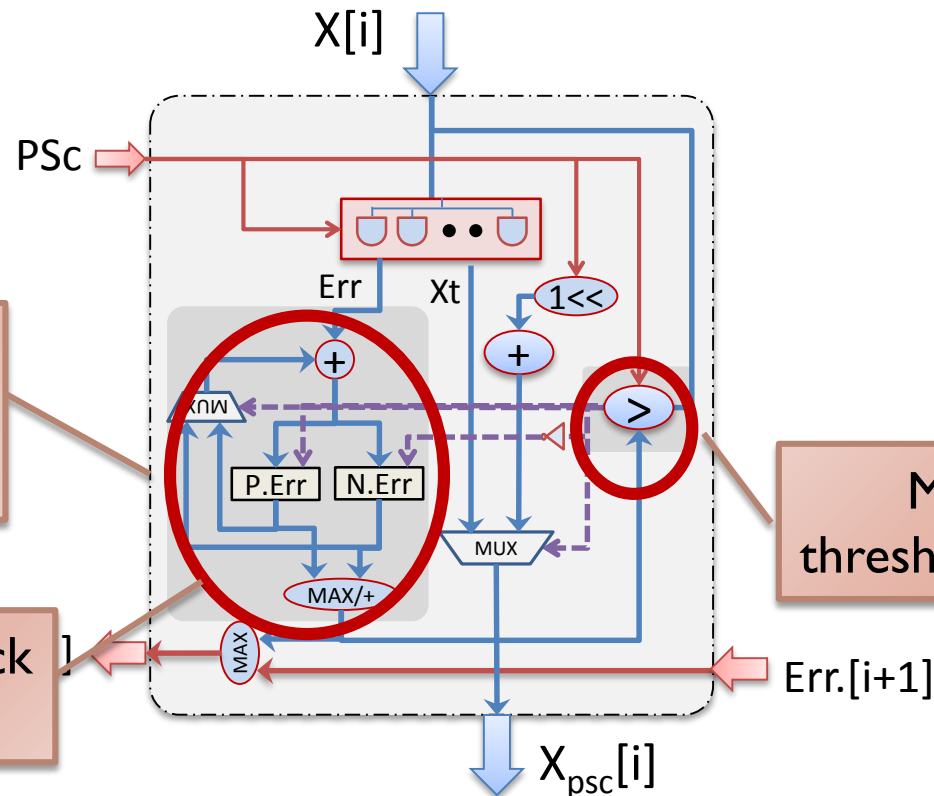
e.g.,  $PSc == 3$

$X[2:0] \geq 4$  ?

Round up : Round down

# PRECISION SCALING WITH ERROR COMPENSATION

- ▶ Quality specified at the output of vector operations
- ▶ **Key idea:** Compensate errors across many scalar operations to reduce overall instruction level error



Track the error in positive and negative directions

Modulate the threshold for round-off

Enables error feedback to software





# QUALITY CONTROL UNIT

- ▶ Set PSc based on required instruction level quality bound

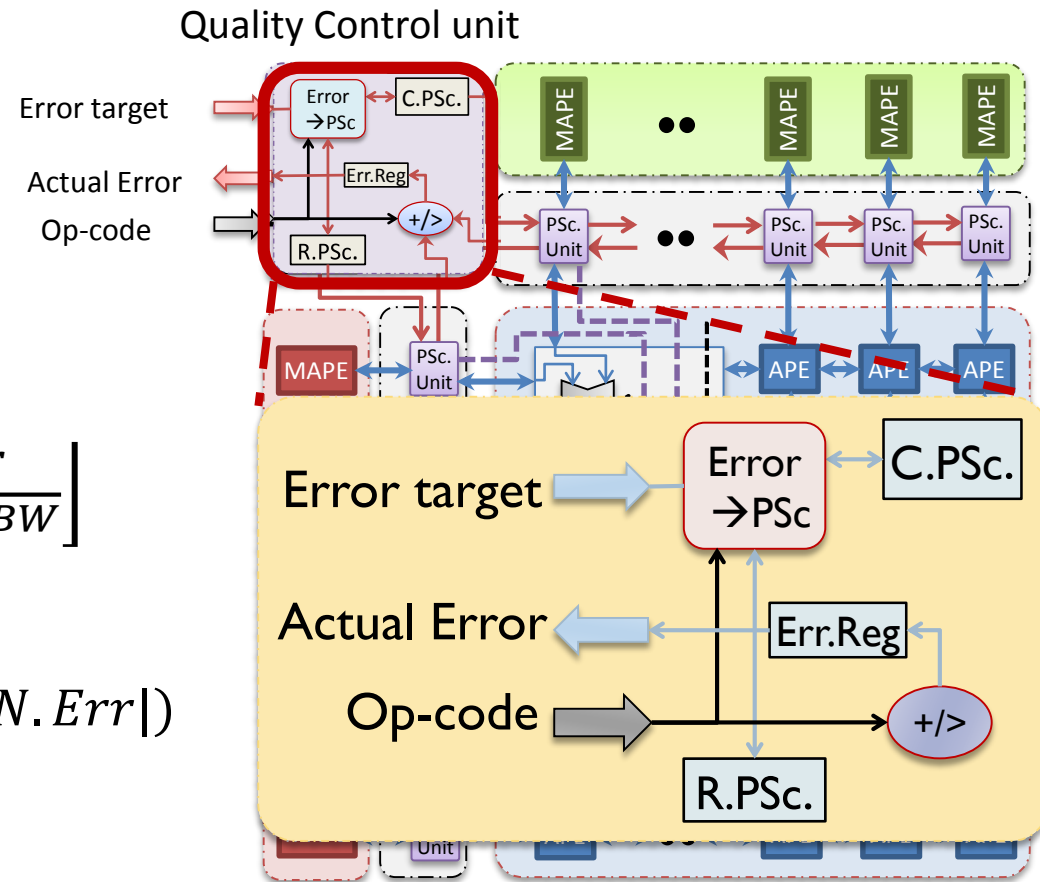
- ▶ E.g., MAC operation

$$\mathbf{C.PSc} = 1 + \lg_2 \left\lceil \frac{MaxErr}{length * 2^{BW}} \right\rceil$$

$$\mathbf{R.PSc} = 0$$

$$\mathbf{E.Reg} = 2^{BW} * \max (|P.Err|, |N.Err|)$$

- ▶ More details in paper



# EXPERIMENTAL METHODOLOGY

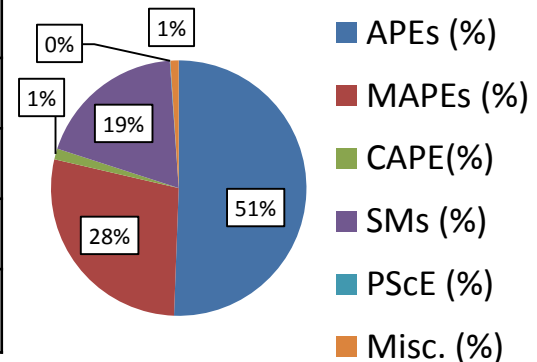
▶ **RTL implementation**  
of QUORA using  
Verilog HDL

▶ Synthesized to **IBM**  
**45nm** technology  
node

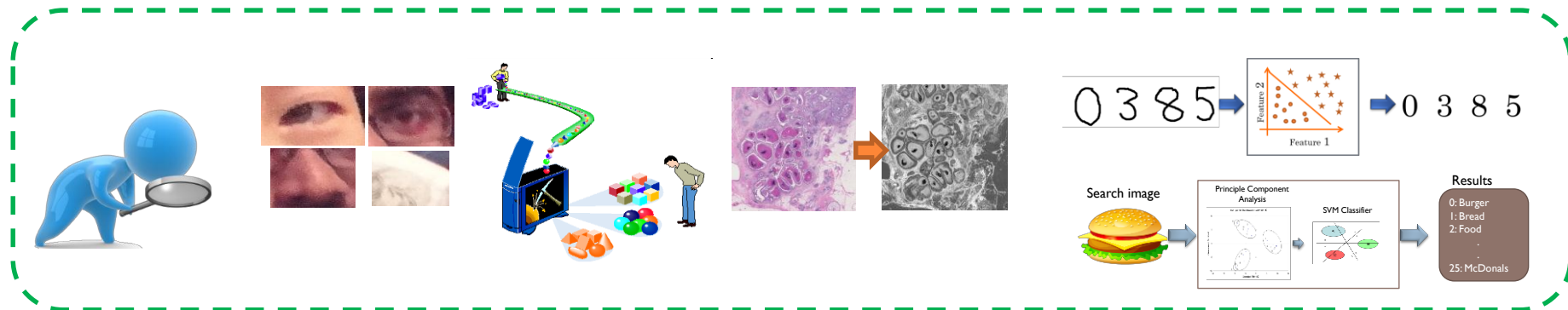
- Design flow: *Synopsys Design Compiler, ModelSim, Synopsys Power Compiler*

Micro-architectural Parameters	Value
Array Dimensions	16 X 16
No. of PEs (APEs + MAPEs+ CAPE)	289 (256 + 32 + 1)
Size of Register File – CAPE / MAPE	32 / 8
No. of SM elements	32
Depth of SM elements	64
Operating Frequency	250 MHz

Metric	Value
Feature Size	45nm
Area	2.6 mm <sup>2</sup>
Power	367.8 mW
Gate Count	502042



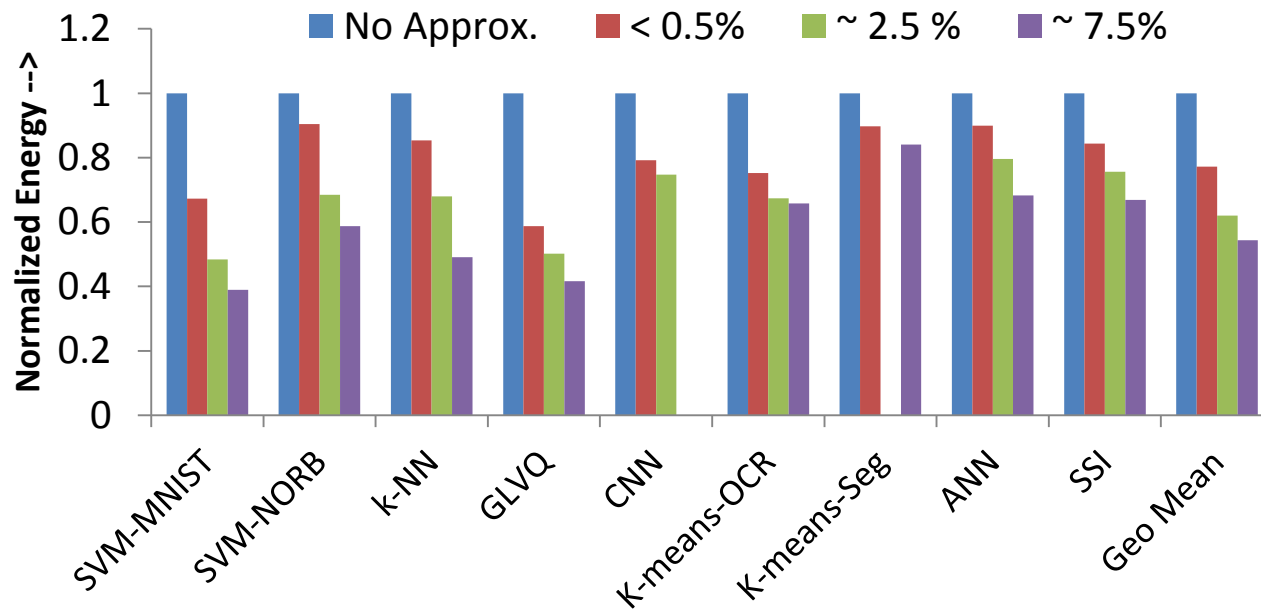
# BENCHMARKS



Applications	Algorithm	Dataset	Quality Metric
Handwritten Digit Recognition (SVM-MNIST)	SVM	MNIST	Percentage classification accuracy
Object Recognition (SVM-NORB)	SVM	NORB	
Digit Classification (CNN)	CNN	MNIST	
Eye Detection (GLVQ)	GLVQ	Image set from NEC labs.	
Optical Character Recognition (k-NN)	K-NN	OCR digits	
Census Data Analysis (ANN)	ANN	Adult	
Document Search (SSI)	SSI	Subset of Wikipedia	No. correct in top 25 results
Image Segmentation (K-Means-Seg)	K-means	Berkeley dataset	Mean distance of clustered points from respective centroids
Optical Character Clustering (K-Means-OCR)	K-means	OCR digits	

# RESULTS SUMMARY

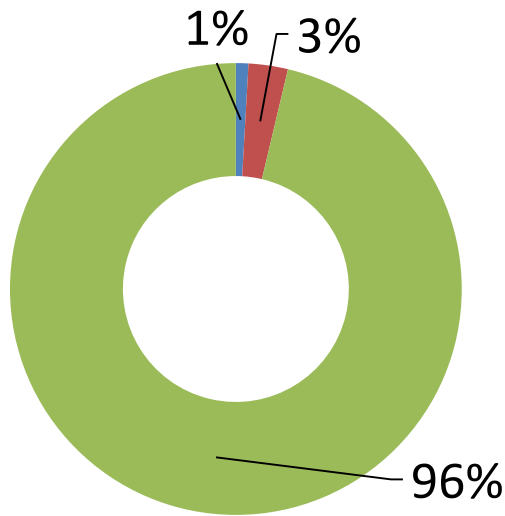
## Energy savings



- ▶ **1.05-1.7X** savings for **NO** loss in output quality
- ▶ **1.18-2.1X** savings for **< 2.5%** quality loss
- ▶ **> 2.5X** savings for **< 7.5%** quality loss

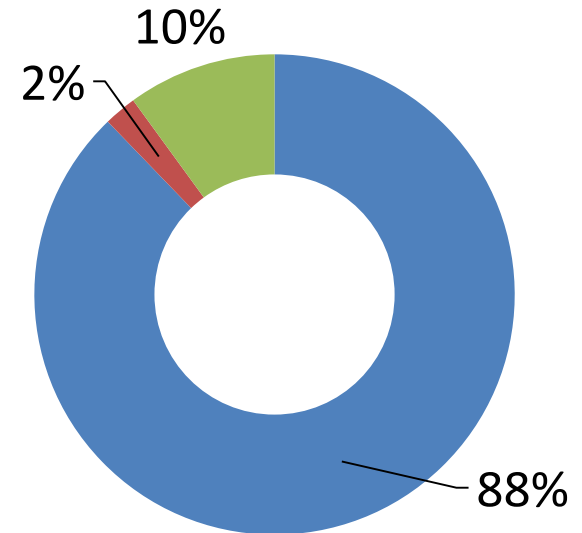
# RESULTS SUMMARY

## QP-instructions in QUORA



■ QP-APE ■ QP-MAPE ■ Accurate

**Dynamic instruction count**



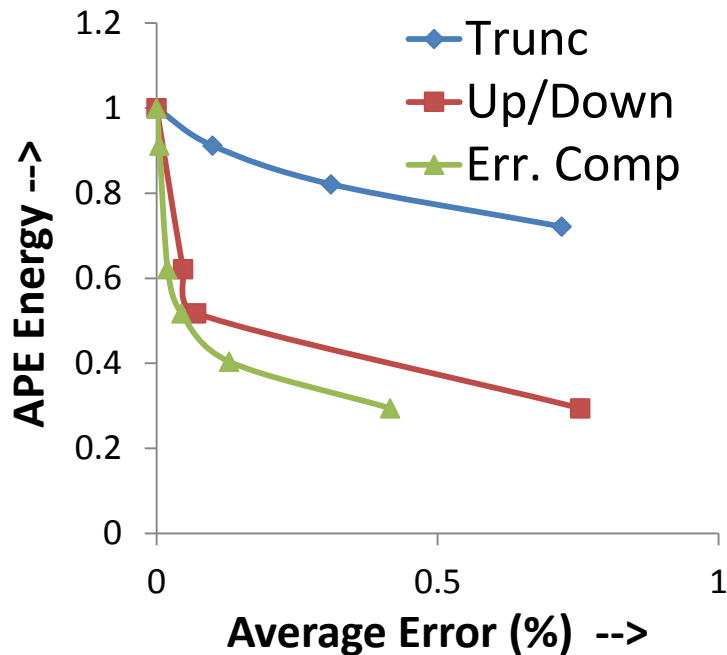
■ QP-APE ■ QP-MAPE ■ Accurate

**Energy**

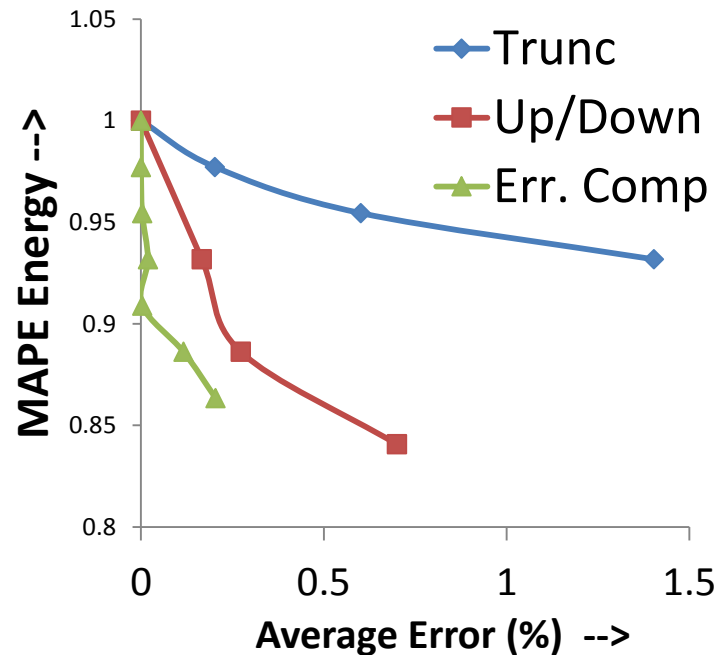
- ▶ **90%** of application energy in quality programmable instructions

# RESULTS SUMMARY

## Precision scaling mechanisms



**MAC - APE**



**ACC - MAPE**

- Precision scaling with error compensation provides superior energy vs. quality trade-off

# SUMMARY

- ▶ **Intrinsic application resilience:** A new dimension to optimize HW and SW
- ▶ **Objective:** Energy-efficient & programmable processor for approximate computing
- ▶ **Quality programmable processors:** Quality codified as part of the instruction set
- ▶ **Quora:** Quality programmable 1D/2D vector processor
  - Quality programmable ISA and microarchitecture